

Constraint Programming for Itemset Mining

Luc De Raedt - Tias Guns – Siegfried Nijssen

Itemset Mining

Input
Items: x, y, z
Transactions:

	x	y	z
{x,y}	1	1	0
{x,y,z}	1	1	1
{x,y}	1	1	0
{y,z}	0	1	1

Search space
lattice wrt. subset relation:

Output
Frequent (min. 3) Itemsets:

{a}	freq=3
{a,d}	freq=3
{d}	freq=4

level-wise enumeration:
+ anti-monotone constraints
X others...

Constraint-based Itemset Mining

Challenge: push constraints into the mining process

- Itemset constraints (minimum size, average cost,...)
- Condensed representations (closed, maximal itemsets,...)
- Alternatives for support (emerging, correlated patterns,...)

...many many, many constraints and algorithms proposed...
BUT **no generic framework**

CP is a general approach to combinatorial problem solving

Input
a model: variables & constraints
 $x_1, x_2, x_3 \in \{1..4\}$
 $x_1 \neq 2$
 $x_1 == x_2$
 $x_2 < x_3$

Output
all valid assignments:
 $\{x_1, x_2, x_3\} = \{1, 1, 2\}, \{1, 1, 3\}, \{3, 3, 4\}, \dots$

Search space
example 4-queens:

constraint-search:

- fix a variable
- propagate constraints
 - if failure: backtrack
 - else: constraint-search

Propagation

- a variable has a domain
- a constraint has a propagator

a propagator can **reduce** the domain of its variables.

Power of CP:

combined effect of propagators
modelling language
 declarative and very expressive

Constraint Programming

FIM_CP: Constraint Programming as Generic Framework for IM

Algorithm 1 Fim_cp's frequent itemset mining model, in Essence'
1: **given** NrT, NrI : int
2: **given** TDB : matrix indexed by [int(1..NrT),int(1..NrI)] of bool
3: **given** Freq : int
4: **find** Items : matrix indexed by [int(1..NrI)] of bool
5: **find** Trans : matrix indexed by [int(1..NrT)] of bool
6: **such that**
7: \$ encode TDB: every Trans its complement has no supported Items
8: **forall** t: int(1..NrT).
9: Trans[t] <=> ((sum i: int(1..NrI). !TDB[t,i]*Items[i]) <= 0),
10: \$ frequency: every Item is supported by sufficiently many Trans
11: **forall** i: int(1..NrI).
12: Items[i] => ((sum t: int(1..NrT). TDB[t,i]*Trans[t]) >= Freq)
Essence' is a known solver-independant constraint modelling language.

$$(8-9) : \forall t \in \mathcal{T} : T_t = 1 \iff \sum_{i \in \mathcal{I}} I_i (1 - D_{ti}) = 0.$$

Propagation:
if one item=1 in transaction's complement \rightarrow trans=0
if all items=0 in transaction's complement \rightarrow trans=1

$$(10-11) : \forall i \in \mathcal{I} : I_i = 1 \rightarrow \sum_{t \in \mathcal{T}} T_t D_{ti} \geq \theta.$$

Propagation:
if too many of item's trans=0 \rightarrow item=0

Many more constraints possible!
(average cost, delta-closed, ...)
Arbitrary combinations of them
(also diff. interpretations)

- Wide range of constraints out-of-the-box
- The modelling power of CP for IM
eg. combining constraints
- Also good on dense datasets
 - Not so good on very large datasets

Conclusion

- Even more constraints
- Optimise CP algorithms for IM
- Other types of patterns
eg. correlated itemsets or trees, graphs, ...

Future Work

